new U.S. Appln. based on Appln. Nos. 60/236,574
& 60/299,360 for Online Intelligent Information...
Express Mail Label No. EL864389877US
Attorney Docket No.: 2102680-990101
Gray Cary et al. – GTS/415-836-2500
Applicants: Victor Hsieh          Sheet 1 of 36

{7,"<B>,""</B>,""<I>,""</I>,""<TABLE>"}

procedure **pcwrap**HLRT (page $P$)
    *scan* to the line 7 in $P$
    while the next occurrence of $li$ in $P$ occurs before the next occurrence of $t$
        scan in $P$ to next occurrence of $l_i$; save position as start of item attribute
        scan in $P$ to next occurrence of $r_i$, save position as end of item attribute
        scan in $P$ to next occurrence of $l_p$; save position as start of price attribute
        scan in $P$ to next occurrence of $r_p$; save position as end of price attribute
    return extracted . . . . ., (item,price), . . . . .} pairs

# FIGURE 1
## (PRIOR ART)

FIGURE 2

new U.S. Appln. based on Appln. Nos. 60/236,574
& 60/299,360 for Online Intelligent Information...
Express Mail Label No. EL864389877US
Attorney Docket No.: 2102680-990101
Gray Cary et al. – GTS/415-836-2500
Applicants: Victor Hsieh        Sheet 3 of 36

100

START

110

RETRIEVE TRAINING
DATA [OFFLINE
DATABASE]

120

MORE
VENDORS TO
LEARN?

NO → END

YES

130

RETRIEVE TRAINING PAGES
[FROM VENDOR SITE]

140

PERFORM INDUCTIVE
LEARNING ON RETRIEVED
TRAINING PAGE

150

STORE LEARNED
RESULT IN OFFLINE
DATABASE

**FIGURE 3**

new U.S. Appln. based on Appln. Nos. 60/236,574
& 60/299,360 for Online Intelligent Information...
Express Mail Label No. EL864389877US
Attorney Docket No.: 2102680-990101
Gray Cary et al. – GTS/415-836-2500
Applicants· Victor Hsieh          Sheet 4 of 36

| Element Name | Element Description |
|---|---|
| Vendor Name | Name of the online shop |
| Vendor URL | The URL of the online shop |
| Form URL | The URL to submit search data |
| Learning Domain | The domain used to learn the vendor description |
| Head | The end of header position of vendor's product pages. This element to learnt to reduce the product page searching space and thus shorten the product information extraction time |
| Tail | The start of footer position of vendor's product pages. Same as Head element, with this element the product information extraction time for Shopper Agent is shorten. |
| Left Delimiter of Item | The Shopper Agent uses this two delimiters to identify a product |
| Right Delimiter of Item | |
| Left Delimiter of Price | The Shopper Agent uses this two delimiters to locate a product's price information |
| Right Delimiter of Price | |

# FIGURE 4

| EXAMPLE | |
|---|---|
| Element Name | Element Description |
| Vendor Name | 800.com |
| Vendor URL | http://www.800.com |
| Form URL | http://www.800.com/search/srchrslts.asp?qs=1&slteentry=All&entry= |
| Learning Domain | Md |
| Head | 1230 |
| Tail | </BODY> |
| Left Delimiter of Item | /td width="12"> |
| Right Delimiter of Item | </fon |
| Left Delimiter of Price | Your Price: $ |
| Right Delimiter of Price | </b> |

# FIGURE 5

new U.S. Appln. based on Appln. Nos. 60/236,574
& 60/299,360 for Online Intelligent Information...
Express Mail Label No. EL864389877US
Attorney Docket No.: 2102680-990101
Gray Cary et al. – GTS/415-836-2500
Applicants: Victor Hsieh          Sheet 5 of 36

140

```
                        ┌──────────────┐
                        │    START     │
                        └──────┬───────┘
                               │
                               ▼                    210
                   ┌───────────────────────┐
                   │ LABEL [POSITION VALUES]│
                   │    TRAINING PAGES      │
                   └───────────┬───────────┘
                               │
                               ▼                    220
                   ┌───────────────────────┐
                   │ GENERATE SET POSSIBLE  │
                   │  VENDOR DESCRIPTION    │
                   │     CANDIDATES         │
                   └───────────┬───────────┘
                               │
                               ▼          230
                          ╱─────────╲
                   NO    ╱   MORE     ╲
              ◄─────────╱   VENDOR     ╲
                        ╲ DESCRIPTION  ╱
                        ╲ CANDIDATES? ╱
                          ╲─────────╱
                               │ YES
                               ▼          240
                   ┌───────────────────────┐
                   │   VALIDATE VENDOR      │
                   │ DESCRIPTION CANDIDATE  │
                   └───────────┬───────────┘
                               │
                               ▼          250
                          ╱─────────╲
                         ╱  VENDOR    ╲
                        ╱ DESCRIPTION  ╲    NO
                        ╲  CANDIDATES  ╱──────►
                        ╲SATISFACTORY?╱
                          ╲─────────╱
                               │ YES
                               ▼
                        ┌──────────────┐
                        │    STOP      │
                        └──────────────┘
```

**FIGURE 6**

new U.S. Appln. based on Appln. Nos. 60/236,574
& 60/299,360 for Online Intelligent Information...
Express Mail Label No. EL864389877US
Attorney Docket No.: 2102680-990101
Gray Cary et al. – GTS/415-836-2500
Applicants: Victor Hsieh          Sheet 6 of 36

| MD PRICE | | `<HTML>` |
| --- | --- | --- |
| | | `<TITLE><B>A Simple Product Catalogs</B></TITLE>` |
| Model Number | PRICE (US$) | `<BODY>` `<H2> MD Price </H2>` |
| HM381MD | 399.95 | `<TABLE BORDER=1>` `<TR BGCOLOR=ORANGE><TH> Model Number </TH>` |
| MD2070 | 599.95 | `<TH> PRICE(US$)</TH></TR>` |
| MD203 | 249.95 | `<TR><TD><B> HM381MD</B></TD><TD><I>399.95</I></TD></TR>` |
| MDR3 | 399.95 | `<TR><TD><B> MD2070</B></TD><TD><I>599.95</I></TD></TR>` |
| | | `<TR><TD><B> MD203</B></TD><TD><I>249.95</I></TD></TR>` |
| | | `<TR><TD><B> MDR3</B></TD><TD><I>399.95</I></TD></TR>` |
| | | `</TABLE>` `<P>` `<HR WIDTH=200 ALIGN=LEFT>` `<P>` `<B> End Of The Product Catalog </B>` `</BODY>` `</HTML>` |

## FIGURE 7

| Product Entry | Labels |
| --- | --- |
| { HM381MD, 399.95} | {<<174, 180>, <197, 202>>} |
| { MD2070, 599.95} | {<<229, 234>, <251, 256>>} |
| { MD203, 249.95} | {<<283, 287>, <304, 309>>} |
| { MDR3, 399.95} | {<<336, 339>, <356, 361>>} |

## FIGURE 8

| Product Entry | Labels |
| --- | --- |
| { PRODUCT, PRICE} | {<PRODUCT LEFT DELIMITER, PRODUCT RIGHT DELIMITER>, <PRICE LEFT DELIMITER, PRICE RIGHT DELIMITER>} |
| * * * | * * * |

## FIGURE 9

new U.S. Appln. based on Appln. Nos. 60/236,574
& 60/299,360 for Online Intelligent Information..
Express Mail Label No. EL864389877US
Attorney Docket No.: 2102680-990101
Gray Cary et al. – GTS/415-836-2500
Applicants: Victor Hsieh          Sheet 7 of 36

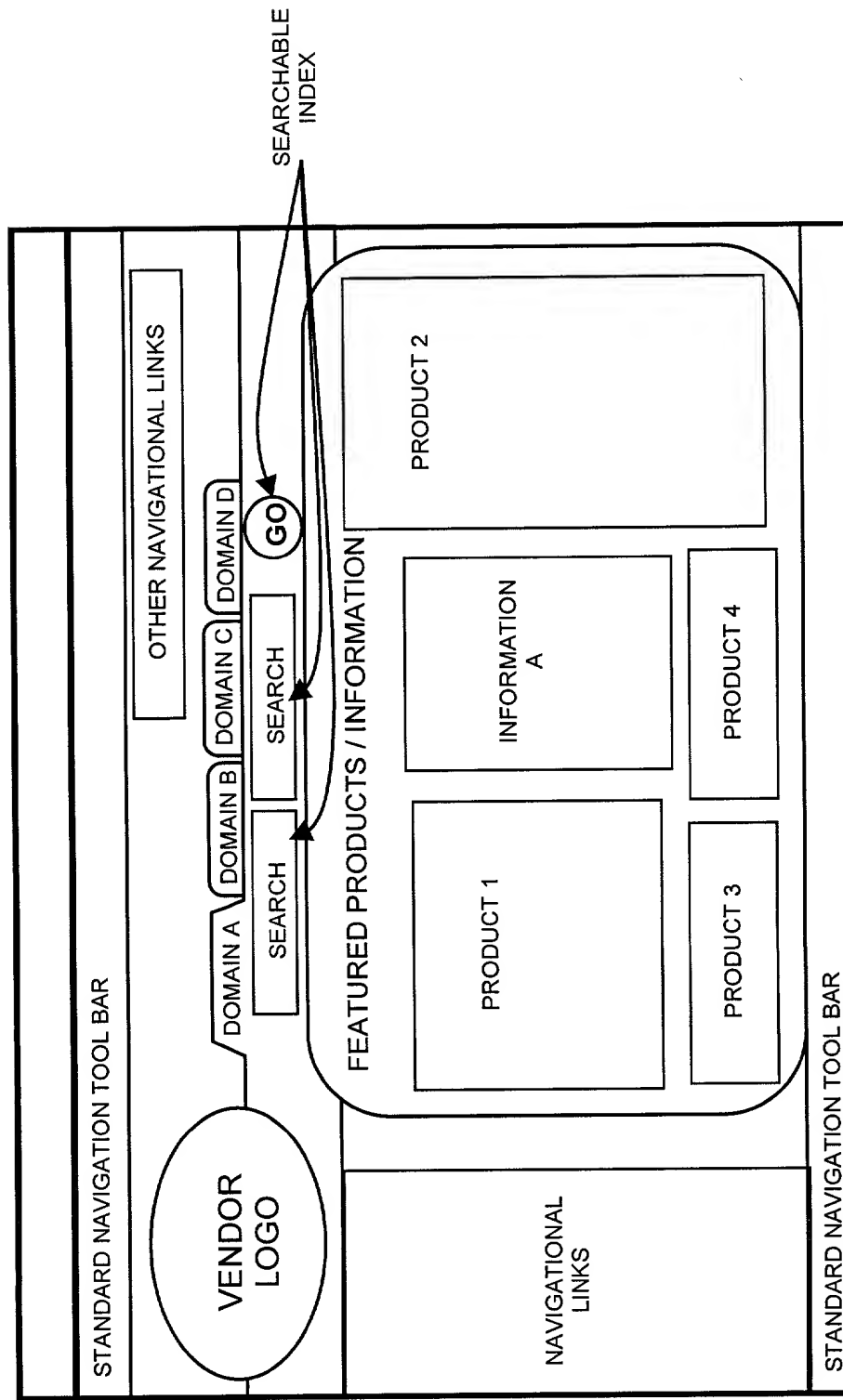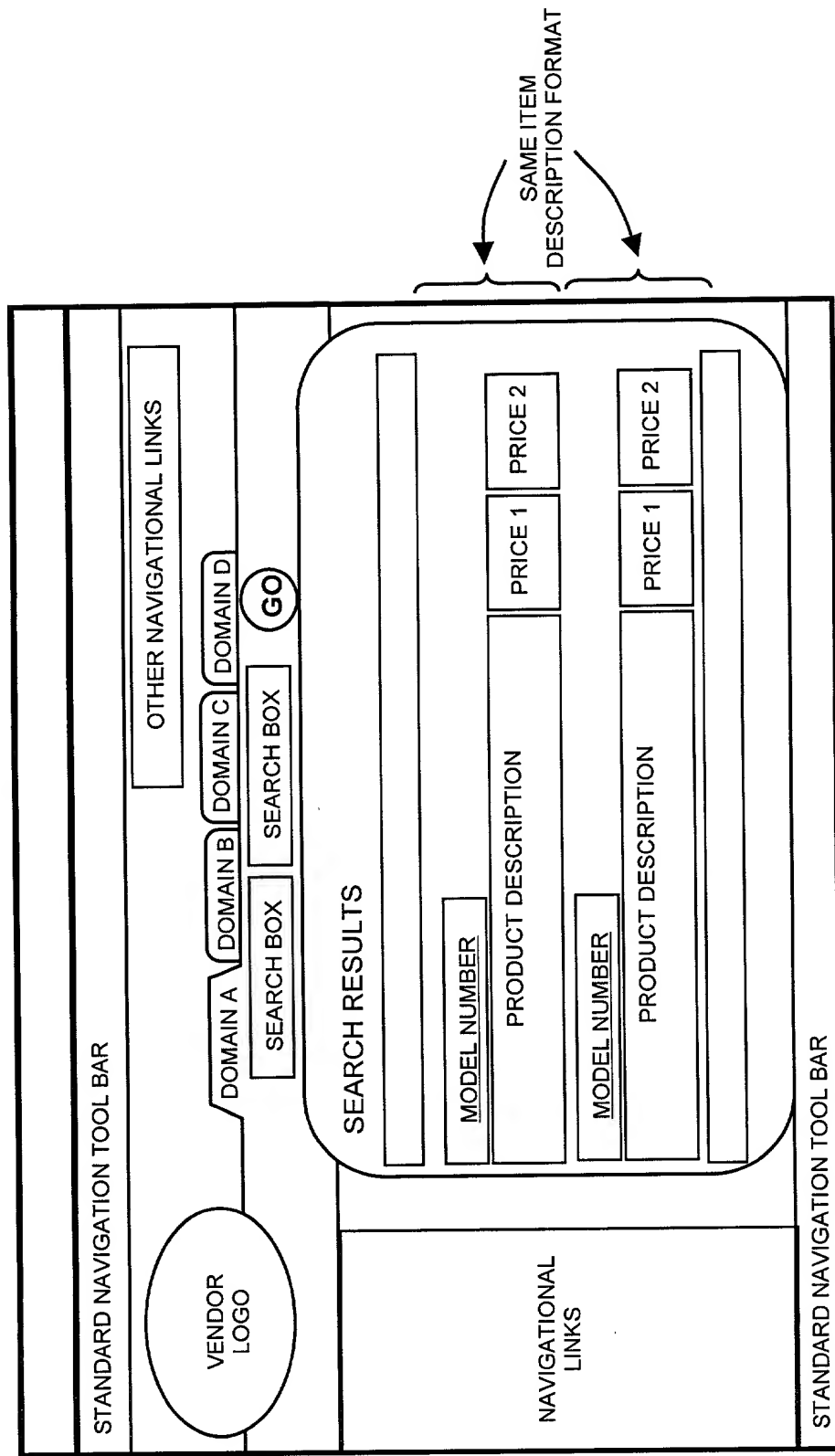| Delimiter | Candidates |
|---|---|
| Item Left Delimiter | <TR><TD><B><br>TR><TD><B><br>R><TD><B><br>><TD><B><br><TD><B><br>TD><B><br>D><B><br>><B><br><B><br>B><br>> |
| Item Right Delimiter | </B></TD><TD><I><br></B></TD><TD><I<br></B></TD><TD><<br></B></TD><TD><br></B></TD><TD<br></B></TD><T<br></B></TD><<br></B></TD><br></B></TD<br></B></T<br></B></<br></B><<br></B><br></B<br></<br>< |
| Price Left Delimiter | </B></TD><TD><I><br>/B></TD><TD><I><br>B></TD><TD><I><br>></TD><TD><I><br></TD><TD><I><br>/TD><TD><I><br>TD><TD><I><br>D><TD><I><br>><TD><I><br><TD><I><br>TD><I><br>D><I><br>><I><br><I><br>I><br>> |

**FIGURE 10A**

new U.S. Appln. based on Appln. Nos. 60/236,574
& 60/299,360 for Online Intelligent Information...
Express Mail Label No EL864389877US
Attorney Docket No.: 2102680-990101
Gray Cary et al. – GTS/415-836-2500
Applicants: Victor Hsieh          Sheet 8 of 36

| Delimiter | Candidates |
|---|---|
| Price Right Delimiter | </I></TD></TR<br></I></TD></T<br></I></TD></<br></I></TD><<br></I></TD><br></I></TD<br></I></T<br></I></<br></I><<br></I><br></I<br></<br>< |
| Head | <<br><H<br><HT<br><HTM<br><HTML<br><HTML><br>... ...<br>... ... ... ... |
| Tail | ><br>L><br>ML><br>TML><br>HTML><br>/HTML><br>... ...<br>... ... ... ... |

## FIGURE 10B

FIGURE 11

STANDARD NAVIGATION TOOL BAR

VENDOR LOGO

OTHER NAVIGATIONAL LINKS

DOMAIN A DOMAIN B DOMAIN C DOMAIN D

SEARCH BOX SEARCH BOX GO

SEARCH RESULTS

MODEL NUMBER

PRODUCT DESCRIPTION

PRICE 1 PRICE 2

MODEL NUMBER

PRODUCT DESCRIPTION

PRICE 1 PRICE 2

NAVIGATIONAL LINKS

STANDARD NAVIGATION TOOL BAR

SAME ITEM DESCRIPTION FORMAT

**FIGURE 12**

STANDARD NAVIGATION TOOL BAR

VENDOR LOGO

OTHER NAVIGATIONAL LINKS

DOMAIN A  DOMAIN B  DOMAIN C  DOMAIN D

SEARCH BOX   SEARCH BOX   GO

SEARCH RESULTS

MODEL NUMBER

PRODUCT DESCRIPTION

PRICE 1   PRICE 2

MODEL NUMBER

PRODUCT DESCRIPTION

PRICE 1   PRICE 2

NAVIGATIONAL LINKS

STANDARD NAVIGATION TOOL BAR

HEAD

CONTENT

TAIL

**FIGURE 13**

## FIGURE 14

DOMAIN DESCRIPTIONS
HM381MD, MD2070

URL
www.800.com

STEP 1

STEP 7

VENDOR DESCRIPTIONS

AGENT/
SYSTEM

STEPS 2, 6

WORLD WIDE WEB
/ INTERNET

STEP 3

Electronics. And more.

STEP 4

STEP 5

```
...
<TR BGCOLOR=ORANGE><TH>Model Number</TH>
<TH>PRICE(US$)</TH></TR>
<TR><TD><B>HM381MD</B></TD><TD><I>399.95</I
></TD></TR>
```

RETURNED PAGE

Applicants: Victor Hsieh        Sheet 12 of 36
Gray Cary et al. – GTS/415-836-2500
Attorney Docket No.: 2102680-990101
Express Mail Label No. EL864389877US
& 60/299,360 for Online Intelligent Information...
new U.S. Appln. based on Appln. Nos. 60/236,574

# Search Results

**Products**
Your search for md returned 19 items  Showing items 1 - 19.

| Product Search Results | Price | |
|---|---|---|
| **Kenwood MASK MD Music System  HM381MD**<br>Kenwood MASK Music System, Self-Hiding Revolving Faceplate, 10 Watts Per Channel, MD Recorder, CD Player, Sophisticated Styling, Silver with Gray Accents, 2-Way High Gloss Speakers, Remote Control<br>*Usually ships within 24 hours* | List Price $500.00<br>Your Price $399.95 | wish list buy now |
| **Kenwood MiniDisc Recorder MD2070**<br>Kenwood MiniDisc Recorder, 10 Second Digital Anti Skip, 24-Bit Resolution A/D and D/A Converters, CD Text Transfer (via Direct Digital Connection), Sampling Rate Converter, Disc and Track Naming, Silver Finish, Remote Control with Jog Shuttle<br>*Usually ships within 24 hours* | List Price $700.00<br>Your Price $599.95 | wish list buy now |
| **Kenwood MiniDisc Recorder MD203**<br>Kenwood MiniDisc Recorder, 10 Second Digital Anti Skip, 20-Bit Resolution A/D and D/A Converters, Sampling Rate Converter, Disc and Track Naming, Black Finish, Remote Control<br>*Usually ships within 24 hours* | List Price $300.00<br>Your Price $249.95 | wish list buy now |

# FIGURE 15A

# FIGURE 15B

## INQUIRY FROM ONLINE HUMAN BUYER/USER

MD or any domains in native characters of buyer/user language



VENDOR DESCRIPTIONS

STEP 1

STEP 2

STEP 8

AGENT/ SYSTEM

STEPS 3,7

WORLD WIDE WEB / INTERNET

STEP 3

STEP 5

STEP 4

```
...
<TR BGCOLOR=ORANGE><TH>Model Number</TH>
<TH>PRICE(US$)</TH></TR>
<TR><TD><B>HM381MD</B></TD><TD><I>399.95</I
></TD></TR>
```

RETURNED PAGE

new U.S. Appln. based on Appln. Nos. 60/236,574
& 60/299,360 for Online Intelligent Information...
Express Mail Label No. EL864389877US
Attorney Docket No.: 2102680-990101
Gray Cary et al. – GTS/415-836-2500
Applicants: Victor Hsieh          Sheet 15 of 36

START

310

300

INITIALIZE CUSTOMER REQUEST

312

ANY HITS IN MEMORY OR CACHE? — YES

NO

320

GET VENDOR DESCRIPTION FROM OFFLINE DATABASE

330

COMPOSE REQUEST WITH VENDOR DESCRIPTION AND CUSTOMER REQUEST

340

FILL OUT REQUEST FORM AT VENDOR SITE AND SUBMIT

360

EXTRACT TARGET INFORMATION FROM RECEIVED DATA

350

DATA RECEIVED WITHIN TIMEOUT?

STORE DATA IN CACHE

358

NO

YES

370

SORT EXTRACTED TARGET INFORMATION

380

GENERATE RESULT PAGES [HTML]

390

DISPLAY RESULT PAGES

STOP

**FIGURE 15C**

new U.S. Appln. based on Appln. Nos. 60/236,574
& 60/299,360 for Online Intelligent Information...
Express Mail Label No. EL864389877US
Attorney Docket No.: 2102680-990101
Gray Cary et al. – GTS/415-836-2500
Applicants: Victor Hsieh          Sheet 16 of 36

**FIGURE 16**

new U.S. Appln. based on Appln. Nos. 60/236,574
& 60/299,360 for Online Intelligent Information...
Express Mail Label No. EL864389877US
Attorney Docket No.: 2102680-990101
Gray Cary et al. -- GTS/415-836-2500
Applicants: Victor Hsieh          Sheet 17 of 36



**Vendor Information**

Vendor Information    Add Vendor

Add Vendor

Details

Vendor Name:    1cache.com

URL:    http://www.1cache.com

Form's URL:    i-bin/nsearch?catalog=1cache&query=

Learning Domain:    dvd

☑ Provide Training Examples

ebook Theater with I-Glasses

GD Dolby Digital DVD Player

bination LD/DVD/CD Player

☐ Enter Wrapper For This Vendor

Head                          Tail

Left Delimiter of Item:        Right Delimiter of Item:

Left Delimiter of Price:       Right Delimiter of Price:

Save

OK        Cancel        Apply

**FIGURE 17**

new U.S. Appln. based on Appln. Nos. 60/236,574
& 60/299,360 for Online Intelligent Information...
Express Mail Label No. EL864389877US
Attorney Docket No.: 2102680-990101
Gray Cary et al. -- GTS/415-836-2500
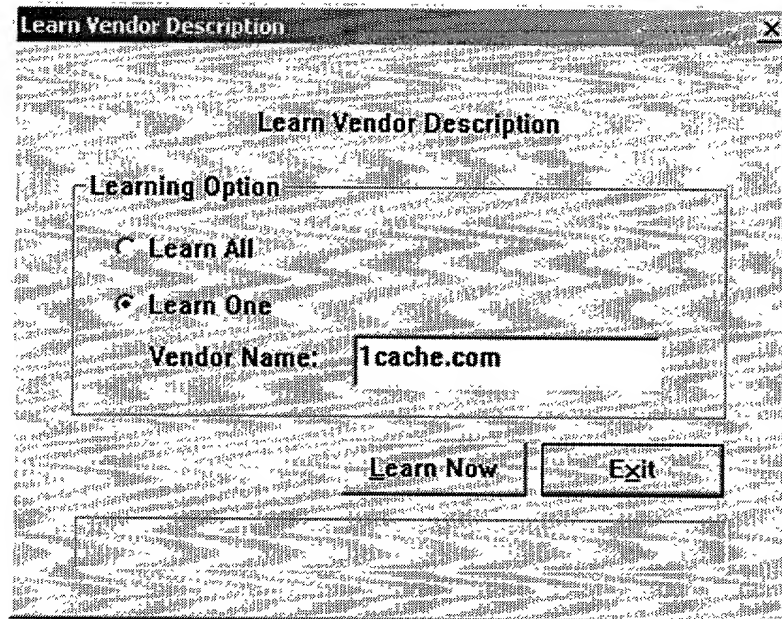Applicants: Victor Hsieh          Sheet 18 of 36

**Vendor Information**

Vendor Information | Add Vendor

Vendor Information

Vendor Name:     1cache.com                          [ Search ]

Details

Vendor Name:       1cache.com

URL:               http://www.1cache.com

Form's URL:        http://st4.yahoo.com/cgi-bin/nsearch?catalog=1ca

Learning Domain:   dvd

Wrapper

Head:                                    Tail:

Left Delimiter of Item:                  Right Delimiter of Item:

Left Delimiter of Price:                 Right Delimiter of Price:

                                                    [ Save ]
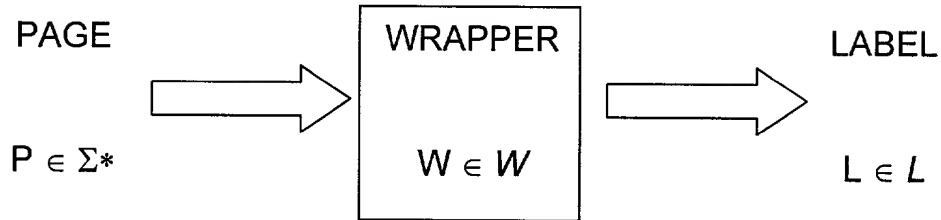
                        [ OK ]    [ Cancel ]    [ Apply ]

**FIGURE 18**

new U.S. Appln. based on Appln. Nos. 60/236,574
& 60/299,360 for Online Intelligent Information...
Express Mail Label No. EL864389877US
Attorney Docket No.: 2102680-990101
Gray Cary et al. – GTS/415-836-2500
Applicants: Victor Hsieh          Sheet 19 of 36

**FIGURE 19**

new U.S. Appln. based on Appln. Nos. 60/236,574
& 60/299,360 for Online Intelligent Information...
Express Mail Label No. EL864389877US
Attorney Docket No.: 2102680-990101
Gray Cary et al. – GTS/415-836-2500
Applicants: Victor Hsieh          Sheet 20 of 36

## Training Data

### Training Example

**Vendor:** 1cache.com                    [ Go ]

| Vendor | Training Example |
|---|---|
| 1cache.com | DVD Virtual Notebook Theater with I-Glas |
| 1cache.com | JVC XV523GD Dolby Digital DVD Player |
| 1cache.com | Pioneer DVL-919 - Combination LD/DVD/ |

[ Add ]
[ Delete ]
[ Edit ]
[ Save ]
[ Cancel ]

**Vendor:** [          ]

**Example:** [          ]

[ Exit ]

## FIGURE 20

## Learn Vendor Description

### Learn Vendor Description

**Learning Option**

( ) Learn All

( ) Learn One

**Vendor Name:** [          ]

[ Learn Now ]    [ Exit ]

## FIGURE 21

new U.S. Appln. based on Appln. Nos. 60/236,574
& 60/299,360 for Online Intelligent Information...
Express Mail Label No. EL864389877US
Attorney Docket No.: 2102680-990101
Gray Cary et al. – GTS/415-836-2500
Applicants: Victor Hsieh          Sheet 21 of 36



**FIGURE 22**

new U.S. Appln. based on Appln. Nos. 60/236,574
& 60/299,360 for Online Intelligent Information...
Express Mail Label No. EL864389877US
Attorney Docket No.: 2102680-990101
Gray Cary et al. – GTS/415-836-2500
Applicants: Victor Hsieh          Sheet 22 of 36

**Vendor Information**

Vendor Information | Add Vendor

Vendor Information

Vendor Name:        1cache.com                                    Search

Details

Vendor Name:        1cache.com

URL:                http://www.1cache.com

Form's URL:         http://st4.yahoo.com/cgi-bin/nsearch?catalog=1ca

Learning Domain:    dvd

Wrapper

Head:                                        Tail:

5230                                         D></TD><TD></TD></TR></

Left Delimiter of Item:                      Right Delimiter of Item:

G SRC=/Img/trans_1x1.gif BORDER=0 WIC        </b

Left Delimiter of Price:                     Right Delimiter of Price:

</b></A></TD>    <TD ALIGN=right><FON       </T

                                             Save

                        OK        Cancel       Apply

**FIGURE 23**

new U.S. Appln. based on Appln. Nos. 60/236,574
& 60/299,360 for Online Intelligent Information...
Express Mail Label No. EL864389877US
Attorney Docket No.: 2102680-990101
Gray Cary et al. – GTS/415-836-2500
Applicants: Victor Hsieh          Sheet 23 of 36

PAGE                    WRAPPER                    LABEL

$P \in \Sigma*$          $W \in \mathcal{W}$          $L \in \mathcal{L}$

## FIGURE 24

---

procedure **exec**$_{HLRT}$ (wrapper (h, $l_i$, $r_i$, $l_p$, $r_p$, t), page $P$)
    $m \leftarrow 0$
    *scan* in $P$ to line $h$
    while the next occurrence of $l_i$ in $P$ occurs before the next occurrence of $t$
        $m \leftarrow m + 1$
        scan in $P$ to the next occurrence of $l_i$ in P; save position as $b_{m,i}$
        scan in $P$ the next occurrence of $r_i$ ; save position as $e_{m,i}$
        scan in $P$ to the next occurrence of $l_p$ in P; save position as $b_{m,p}$
       scan in $P$ to next occurrence of $r_p$; save position as $e_{m,p}$
return label $\{\ldots, <<b_{mi}, e_{mi}>,<b_{mp}, e_{mp}>>\}$

## FIGURE 25

---

Procedure **learnHLRT** (examples $\varepsilon$)

1. Generate the candidate sets **Cands**$_r$(l, $\varepsilon$), **Cands**$_l$(l, $\varepsilon$ ), **Cands**$_r$(p, $\varepsilon$), **Cands**$_l$(p, $\varepsilon$ ).

2. Enumerate the cross product of these candidate sets; each element W = (h, $l_i$, $r_i$, $l_p$, $r_p$, t), of this cross product is a wrapper. Halt if W is satisfactory i.e., execLR(W, $P_n$) = Ln, for every ($P_n$, Ln) $\in \varepsilon$.

## FIGURE 26

new U.S. Appln. based on Appln. Nos. 60/236,574
& 60/299,360 for Online Intelligent Information...
Express Mail Label No. EL864389877US
Attorney Docket No.: 2102680-990101
Gray Cary et al. – GTS/415-836-2500
Applicants: Victor Hsieh          Sheet 24 of 36

```
/*
 * The main problem
 */
    procedure learnHLRT (examples ε) for each 1 <k <K
        for each u ∈ cands₁(p, ε): if valid₁(u, p, ε) then lₚ← u and terminate this loop
        for each u ∈ candsᵣ(i, ε): if validᵣ(u, i, ε) then rᵢ←and terminate this loop
        for each u ∈ candsᵣ(p, ε): if validᵣ(u, p, ε) then rₚ←u and terminate this loop
        for each uₗᵢ ∈ cands₁(i, ε)
        for each uₕ ∈ candsₕ(ε)
            for each uₜ ∈ candsₜ(ε)
                if validᵢₕ,ₜ( uₗᵢ, uₕ, uₜ, ε) then
                    lᵢ←uₗᵢ, h←uₕ, t←uₜ, and terminate these three loops
        return HLRT wrapper (h, lᵢ, rᵢ, lₚ, rₚ, t)


/*
 *Generate a set of candidates for left delimiter of price attribute
 */
procedure cands₁(attribute price, examples ε)
 return the set of all suffixes of the shortest string in neighbors₁(price, ε)


/*
 *Generate a set of candidates for right delimiter of attribute a.  Here a could be item or price
 */
procedure candsᵣ(attribute a, examples ε)
     return the set of all prefixes of the shortest string in neighborsᵣ(a, ε)
/*
 * Generate a set of candidates for a page's head
 */
procedure candsₕ (examples ε)
     return the set of all substrings of the shortest string in heads(ε)


/*
 *Generate a set of candidates for a page's tail
 */
procedure candsₜ(example ε)
     return the set of all substrings of the shortest string in tails(ε)


/*
 *Determine whether a particular candidate for the left delimiter of price is valid
 *This procedure applies constraints C3
 */
    procedure valid(candidate u, attribute price, examples ε)
            for each s∈ neighbors₁(price, ε): if u is not a proper suffix of s then return FALSE
```

## Figure 27A

new U.S. Appln. based on Appln. Nos. 60/236,574
& 60/299,360 for Online Intelligent Information...
Express Mail Label No. EL864389877US
Attorney Docket No.: 2102680-990101
Gray Cary et al. – GTS/415-836-2500
Applicants: Victor Hsieh          Sheet 25 of 36

```
    for each s ∈ tails(ε): If u is a sub string of s then return FALSE
    return TRUE

/*
 *Determine whether a particular candidate for the right delimiter of item or price is
 *valid
 *This procedure applies constraints C1 and C2
*/
    procedure valid_r(candidate u, attribute a, examples ε)
for each s ∈ attribs(a, ε): if u is a sub string of s then return FALSE
for each s ∈ neighbors(a, ε): if u is not a prefix of s then return FALSE
    return TRUE
/*
 *Determine whether a particular combination of candidate uh, ut, and uli for h, t, and
 *li respectively are satisfactory.
 */
procedure valid_li,h,t(candidates u_li, u_h, u_t, example ε)
    for each s ∈ heads(ε)
        if is not a sub string of s then return FALSE
        if U_l1 is not a proper suffix of scan (s, Uh) then return FALSE
            if Ut occurs before U1_1 in scan (s, Uh), then return FALSE
    for each s ∈ tails (ε)
        if Ut 1S not a sub string of s then return FALSE
        if Uli occurs before Ut in s then return FALSE
    for each s ∈ seps(ε)
        if Uli is not a proper suffix of s then return FALSE
        if Ut occurs before Uli in s then return FALSE
    return TRUE

/*
 *Return a set of containing all values of a attribute, either item(i) or price(p) in each
 *example
*/.
    procedure attribs(attribute a, examples ε)
        return U_(Pn,Ln) ∈ ε(Pn[b_m,a, e_m,a] I<. .., (b_m,k, e_m,k), . . . > ∈ Ln}

/*
 *Return all strings to the left of an attribute, whether these strings are in the heads
 * or the bodies of the pages
*/
    procedure neighbors_i(attribute a, examples ε)
        if a=i then return seps (i, ε) U heads (ε) else return seps (a, ε)
/*
```

**FIGURE 27B**

new U.S. Appln. based on Appln. Nos. 60/236,574
& 60/299,360 for Online Intelligent Information...
Express Mail Label No. EL864389877US
Attorney Docket No.: 2102680-990101
Gray Cary et al. – GTS/415-836-2500
Applicants: Victor Hsieh          Sheet 26 of 36

NEXTGEN.ASP

BROWSER

PROJECT.ASP

INTERNET
INFORMATION
SERVER

ACTIVE SERVER
PAGES

MSWCRUM.DLL

WEBCLASS
RUNTIME

— APPLICATION OBJECT
— BROWSERTYPE OBJECT
— REQUEST OBJECT
— RESPOND OBJECT
— SERVER OBJECT
— SESSION OBJECT

WEB SERVER

## FIGURE 28

CLIENT

CONNECTABLE
OBJECT

LINK

ICONNECTION POINT

CONNECTION
POINT
OBJECT

OUTGOING INTERFACE

## FIGURE 29

new U.S. Appln. based on Appln. Nos. 60/236,574
& 60/299,360 for Online Intelligent Information...
Express Mail Label No. EL864389877US
Attorney Docket No.: 2102680-990101
Gray Cary et al. – GTS/415-836-2500
Applicants: Victor Hsieh          Sheet 27 of 36

**FIGURE 30**

```
┌─────────┐
│  START  │
└─────────┘
     │
     ▼
┌──────────────────────────────────────┐
│        Call SQLAccEnv                │
│  to allocate environment workspace   │
│     (provides HENV handle)           │
└──────────────────────────────────────┘
     │
     ▼
┌──────────────────────────────────────┐
│       Call SQLAllocConnect           │
│  to allocate connection workspace    │
│     (provides HDBC handle)           │
└──────────────────────────────────────┘
     │
     ▼
┌──────────────────────────────────────┐
│         Call SQLConnect              │
│     to connect to the database       │
└──────────────────────────────────────┘
     │
     ▼
┌──────────────────────────────────────┐
│        Call SQLAllocStmt             │
│  to allocate a statement workspace   │
│      (provides HSTMT handle)         │
└──────────────────────────────────────┘
     │
     ▼
┌──────────────────────────────────────┐
│     Call ODBC member functions       │
│ (SQLExecDirect, SQLFetch, SQLGoData) │
│   to perform a database operation    │
└──────────────────────────────────────┘
     │
     ▼
┌──────────────────────────────────────┐
│        Call SQLFreeStmt              │
│  to deallocate statement workspace   │
│     (invalidates HSTMT handle)       │
└──────────────────────────────────────┘
     │
     ▼
┌──────────────────────────────────────┐
│        Call SQLDisconnect            │
│   to disconnect to the database      │
└──────────────────────────────────────┘
     │
     ▼
┌──────────────────────────────────────┐
│       Call SQLFreeConnect            │
│  to deallocate connection workspace  │
│     (invalidates HDBC handle)        │
└──────────────────────────────────────┘
     │
     ▼
┌──────────────────────────────────────┐
│         Call SQLFreeEnv              │
│ to deallocate environment workspace  │
│     (invalidates HENV handle)        │
└──────────────────────────────────────┘
     │
     ▼
┌─────────┐
│  STOP   │
└─────────┘
```

new U.S. Appln. based on Appln. Nos. 60/236,574
& 60/299,360 for Online Intelligent Information...
Express Mail Label No. EL864389877US
Attorney Docket No.: 2102680-990101
Gray Cary et al. – GTS/415-836-2500
Applicants: Victor Hsieh          Sheet 28 of 36

**FIGURE 31**

```
                    ┌──────────┐
                    │  START   │
                    └──────────┘
                         │
                         ▼
        ┌───────────────────────────────────┐
        │                                   │──┐
        │   Create CInternetSession objects │   ) 1002
        │                                   │──┘
        └───────────────────────────────────┘
                         │
                         ▼
        ┌───────────────────────────────────┐
        │      Call member function         │──┐
        │      s.GetHttpConnection           │   ) 1004
        │  (creates CHttpConnection object c)│──┘
        └───────────────────────────────────┘
                         │
                         ▼
        ┌───────────────────────────────────┐
        │      Call member function         │──┐
        │      c.OpenRequest                 │   ) 1006
        │   (creates CHttpFile objectf)      │──┘
        └───────────────────────────────────┘
                         │
                         ▼
        ┌───────────────────────────────────┐
        │      Call member function         │──┐
        │      f.SendRequest                 │   ) 1008
        │ (sends the POST request and formdata)│──┘
        └───────────────────────────────────┘
                         │
                         ▼
        ┌───────────────────────────────────┐
        │      Call member function         │◄────┐
        │      f.Read                        │──┐  │
        │   (returns chunk of response)      │   )1010
        └───────────────────────────────────┘      │
                         │                          │
                         ▼                          │
                      ◇◇◇◇◇                          │
                   ◇         ◇      No   ┌──────────────────┐
                  ◇ Zero bytes ◇────────►│   Process data   │──) 1014
                   ◇ returned? ◇         └──────────────────┘
                      ◇◇◇◇◇
                        │ 1012
                        │ Yes
                        ▼
                    ┌──────────┐
                    │  STOP    │
                    └──────────┘
```

Choose a Channel

All Ontological Domains ▽

Advanced Search

WEBSITE LOGO

QUICK SEARCH

Enter a "MULTILIGUAL CHARACTER" keyword and click GO!

Keyword: [ ] (GO!)

SEARCH TIPS

• • •

MEMBER SIGN IN

[ ]
[ ]
[Login] [Reset]

Registered Vendor Links
Abcd.com      cdef.com
1234z.com    23ezid..com
pkkw.com     er5lcin.com

WEBSITE INTRODUCTORY
TEXT AND IMAGES

Reverse Advertising Sites

⊙Feedback

[ ]

[Submit]

**FIGURE 32**

# WEBSITE LOGO

NAVIGATION BUTTONS TO OTHER PAGES IN WEBSITE

Choose a Channel

ALL G2B PROCUREMENT DOMAINS ▽

Advanced Search

| DOMAIN A | DOMAIN B | DOMAIN C | DOMAIN D | OTHER DOMAINS |
|---|---|---|---|---|

QUICK SEARCH

Enter a "MULTILIGUAL CHARACTER" keyword and click GO!

Keyword: [ ]  (GO!)

SEARCH TIPS

• • •

Company A
Infomation about G2B Company A

**Company G**
**Information about G2B Company G**

Reverse Advertising Sites

(▲) Registered Vendor Links
Abcd.com    cdef.com
1234z.com   23ezid..com
pkkw.com    er5lcin.com

(▲) Feedback

[ Submit! ]

## FIGURE 33

|  | NAVIGATION BUTTONS TO OTHER PAGES IN WEBSITE |  |
|---|---|---|

**WEBSITE LOGO**

**Choose a Channel**

All G2B Tender Procurement Domains ▼

**Advanced Search**

| DOMAIN A | DOMAIN B | DOMAIN C | DOMAIN D | OTHER DOMAINS |
|---|---|---|---|---|

ADVANCED SEARCH AGENTS ARE ON!

Enter a "MULTILIGUAL CHARACTER" keyword and click GO!

Multilingual Government Networks

| ☐ Select All | | | |
|---|---|---|---|
| Public Body | Speed | Air/Shipping/Ground/Free Overnight Del. | Ease of Navigation |
| ☐ G2B Company G | ⇧⇧⇧⇧ | ⇧ | EEEE |

Time Out

Reserved Price Range    From $            To $

Manufacturer or OEM

Keyword:

( GO! )

SEARCH TIPS

* * *

**FIGURE 34**

# WEBSITE LOGO

NAVIGATION BUTTONS TO OTHER PAGES IN WEBSITE

Choose a Channel

All Ontological Domains ▼

Advanced Search

| DOMAIN A | DOMAIN B | DOMAIN C | DOMAIN D | OTHER DOMAINS |
| --- | --- | --- | --- | --- |

## QUICK SEARCH

Enter a "MULTILIGUAL CHARACTER" keyword and click GO!

Keyword: [____] GO!

SEARCH TIPS

* * *

For B2B Trading Exchanges, Auctions, Vertical Portals, Production and Non-Production e-Procurement Supplies

Agricultural / Farming
Apparel
Automotive/ Transportation
Business Machines
Computers & Electronics
Industrial Equipment
Medical & Scientivic Equipment

More Business to Business
MRO/Raw Materials
Office Furnishings
Office Supplies
Real Estate
Restaurant/Hospitality
Services
Travel & Entertainment

Reverse Advertising Sites

Registered Vendor Links
Abcd.com    cdef.com
1234z.com   23ezid..com
pkkw.com    er5lcin.com

Feedback

Submit!

## FIGURE 35

| | NAVIGATION BUTTONS TO OTHER PAGES IN WEBSITE | | | |

# WEBSITE LOGO

Choose a Channel

| All Ontological Domains | ▽ |

Advanced Search

| DOMAIN A | DOMAIN B | DOMAIN C | DOMAIN D | OTHER DOMAINS |

ADVANCED SEARCH AGENTS ARE ON!

Enter a "MULTILIGUAL CHARACTER" keyword and click GO!

☐ Select All    Multilingual B2B Suppliers

| Online Stores | Speed | Air/Shipping/Ground/Free Overnight Del. | Ease of Navigation |
|---|---|---|---|
| ☐ B2B Company 1 | ⇑⇑⇑⇑⇑ | | EEEE |
| ☐ B2B Company 2 | ⇑⇑⇑⇑⇑⇑ | ⇑ | EEEEE |
| ☐ B2B Company 3 | ⇑⇑⇑⇑⇑ | ⇑ | EEEE |
| ☐ B2B Company 4 | ⇑⇑⇑⇑⇑ | ⇑ | EEEE |

SEARCH TIPS

* * *

Time Out

Price Range    From $    To $

Manufacturer

Keyword:    ( GO! )

# FIGURE 36

**FIGURE 37**

# WEBSITE LOGO

Choose a Channel

All Domains ▽

Advanced Search

| DOMAIN A | DOMAIN B | DOMAIN C | DOMAIN D | OTHER DOMAINS |
| --- | --- | --- | --- | --- |

ADVANCED SEARCH AGENTS ARE ON!

Enter a "MULTILIGUAL CHARACTER" keyword and click GO!

**English Vendors**

| Online Stores | Speed | Air/Shipping/Ground/Free Overnight Del. | Ease of Navigation |
| --- | --- | --- | --- |
| ☐ Select All | | | |
| ☐ Store 1AAC | ʃʃʃʃ | | EEEE |
| ☐ Store 1443A | ʃʃʃʃʃʃ | ⇑ | EE |
| ☐ Store 1443A | ʃʃʃʃʃ | ⇑ | EE |

**Chinese Vendors**

| | | | |
| --- | --- | --- | --- |
| ☐ Select All | | | |
| ☐ Store C5A3 | ʃʃʃʃʃʃ | ⇑ | EEEE |
| ☐ Store DRGWS4 | ʃʃʃʃʃ | ⇑ | EEEE |

Time Out [ ]

Price Range     From $ [ ]     To $ [ ]

Manufacturer [ ]

Keyword: [ ]     ( GO! )

SEARCH TIPS

\* \* \*

## FIGURE 38

| | NAVIGATION BUTTONS TO OTHER PAGES IN WEBSITE | | | |
|---|---|---|---|---|

## WEBSITE LOGO

Choose a Channel

| DOMAIN A | DOMAIN B | DOMAIN C | DOMAIN D | OTHER DOMAINS |
|---|---|---|---|---|

ALL DOMAINS ▼

Advanced Search

QUICK SEARCH

Enter a "MULTILIGUAL CHARACTER" keyword and click GO!

Keyword: [____] (GO!)

SEARCH TIPS

* * *

## Search Results

**Products**
Your search for XX returned 5 Domain C items.

| Model / Features | Store | Price |
|---|---|---|
| AMS7 XX Widget<br>[Features] | Store 1443A | $399.95<br>Buy now |
| AMS2 XX Widget<br>[Features] | Store 1443A | $329.95<br>Buy now |
| AMS7 XX Widget<br>[Features] | Store DRGWS4 | $459.95<br>Buy now |
| AMS7 XX Widget<br>[Features] | Store 4U$U$ | $499.95<br>Buy now |
| ZRTQ XX Widget<br>[Features] | Store OOHE | $529.95<br>Buy now |

## FIGURE 39